

DATABASE ANNOTATION AND RETRIEVAL

The present invention relates to the annotation of data files which are to be stored in a database for facilitating their subsequent retrieval. The present invention is also concerned with a system for generating the annotation data which is added to the data file and to a system for searching the annotation data in the database to retrieve a desired data file in response to a user's input query.

Databases of information are well known and suffer from the problem of how to locate and retrieve the desired information from the database quickly and efficiently. Existing database search tools allow the user to search the database using typed keywords. Whilst this is quick and efficient, this type of searching is not suitable for various kinds of databases, such as video or audio databases.

According to one aspect, the present invention aims to provide a data structure which will allow the annotation of data files within a database which will allow a quick and efficient search to be carried out in response to a user's input query.

According to one aspect, the present invention provides data defining a phoneme and word lattice for use as an annotation data for annotating data files to be stored within a database. Preferably, the data defines a plurality of nodes within the lattice and a plurality of links connecting the nodes within the lattice and further data associates a plurality of phonemes with a respective plurality of links and further data associates at least one word with at least one of said links.

According to another aspect, the present invention provides a method of searching a database comprising the annotation data discussed above, in response to an input query by a user. The method preferably comprises the

5 steps of generating phoneme data and word data corresponding to the user's input query; searching the database using the word data corresponding to the user's query; selecting a portion of the data defining the phoneme and word lattice in the database for further

10 searching in response to the results of the word search; searching said selected portion of the database using said phoneme data corresponding to the user's input query; and outputting the search results.

15 According to this aspect, the present invention also provides an apparatus for searching a database which employs the annotation data discussed above for annotating data files therein. The apparatus preferably comprises means for generating phoneme data and word data

20 corresponding to a user's input query; means for searching the database using the word data corresponding to the user's input query to identify similar words within the database; means for selecting a portion of the annotation data in the database for further searching in

25 response to the results of the word search; means for searching the selected portion using the phoneme data corresponding to the user's input query; and means for outputting the search results.

30 The phoneme and annotation data for a data file may be generated from the data file itself or from a typed or spoken annotation input by the user.

Exemplary embodiments of the present invention will now

35 be described with reference to the accompanying figures,

in which:

Figure 1 is a schematic view of a computer which is programmed to operate an embodiment of the present invention;

Figure 2 is a block diagram showing a phoneme and word annotator unit which is operable to generate phoneme and word annotation data for appendage to a data file;

Figure 3 is a block diagram illustrating one way in which the phoneme and word annotator can generate the annotation data from an input video data file;

Figure 4a is a schematic diagram of a phoneme lattice for an example audio string from the input video data file;

Figure 4b is a schematic diagram of a word and phoneme lattice embodying one aspect of the present invention, for an example audio string from the input video data file;

Figure 5 is a schematic block diagram of a user's terminal which allows the user to retrieve information from the database by a voice query;

Figure 6a is a flow diagram illustrating part of the flow control of the user terminal shown in Figure 5;

Figure 6b is a flow diagram illustrating the remaining part of the flow control of the user terminal shown in Figure 5;

Figure 7 is a flow diagram illustrating the way in which a search engine forming part of the user's terminal

carries out a phoneme search within the database;

Figure 8 is a schematic diagram illustrating the form of a phoneme string and four M-GRAMS generated from the  
5 phoneme string;

Figure 9 is a plot showing two vectors and the angle between the two vectors;

10 Figure 10 is a schematic diagram of a pair of word and phoneme lattices, for example audio strings from two speakers;

Figure 11 is a schematic block diagram illustrating a  
15 user terminal which allows the annotation of a data file with annotation data generated from an audio signal input from a user;

Figure 12 is a schematic diagram of phoneme and word  
20 lattice annotation data which is generated for an example utterance input by the user for annotating a data file;

Figure 13 is a schematic block diagram illustrating a  
25 user terminal which allows the annotation of a data file with annotation data generated from a typed input from a user;

Figure 14 is a schematic diagram of phoneme and word  
30 lattice annotation data which is generated for a typed input by the user for annotating a data file;

Figure 15 is a block schematic diagram showing the form of a document annotation system;

35 Figure 16 is a block schematic diagram of an alternative

document annotation system;

Figure 17 is a block schematic diagram of another document annotation system;

5

Figure 18 is a schematic block diagram illustrating a user terminal which is operable to access a database located on a remote server via a data network in response to an input utterance by the user;

10

Figure 19 is a schematic block diagram of a user terminal which allows a user to access a database located in a remote server in response to an input utterance from the user;

15

Figure 20 is a schematic block diagram of a user terminal which allows a user to access a database by a typed input query; and

20

Figure 21 is a schematic block diagram illustrating the way in which a phoneme and word lattice can be generated from script data contained within a video data file.

Embodiments of the present invention can be implemented using dedicated hardware circuits, but the embodiment to be described is implemented in computer software or code, which is run in conjunction with processing hardware such as a personal computer, work station, photocopier, facsimile machine, personal digital assistant (PDA) or the like.

30

Figure 1 shows a personal computer (PC) 1 which is programmed to operate an embodiment of the present invention. A keyboard 3, a pointing device 5, a microphone 7 and a telephone line 9 are connected to the

35

- PC 1 via an interface 11. The keyboard 3 and pointing device 5 enable the system to be controlled by a user. The microphone 7 converts acoustic speech signals from the user into equivalent electrical signals and supplies them to the PC 1 for processing. An internal modem and speech receiving circuit (not shown) is connected to the telephone line 9 so that the PC 1 can communicate with, for example, a remote computer or with a remote user.
10. The programme instructions which make the PC 1 operate in accordance with the present invention may be supplied for use with an existing PC 1 on, for example, a storage device such as a magnetic disc 13, or by downloading the software from the Internet (not shown) via the internal modem and telephone line 9.

#### DATA FILE ANNOTATION

- Figure 2 is a block diagram illustrating the way in which annotation data 21 for an input data file 23 is generated in this embodiment by a phoneme and word annotating unit 25. As shown, the generated phoneme and word annotation data 21 is then combined with the data file 23 in the data combination unit 27 and the combined data file output thereby is input to the database 29. In this embodiment, the annotation data 21 comprises a combined phoneme (or phoneme like) and word lattice which allows the user to retrieve information from the database by a voice query. As those skilled in the art will appreciate, the data file 23 can be any kind of data file, such as, a video file, an audio file, a multimedia file etc.

- A system has been proposed to generate N-Best word lists for an audio stream as annotation data by passing the audio data from a video data file through an automatic

speech recognition unit. However, such word-based systems suffer from a number of problems. These include (i) that state of the art speech recognition systems still make basic mistakes in recognition; (ii) that state of the art automatic speech recognition systems use a dictionary of perhaps 20,000 to 100,000 words and cannot produce words outside that vocabulary; and (iii) that the production of N-Best lists grows exponentially with the number of hypothesis at each stage, therefore resulting in the annotation data becoming prohibitively large for long utterances.

The first of these problems may not be that significant if the same automatic speech recognition system is used to generate the annotation data and to subsequently retrieve the corresponding data file, since the same decoding error could occur. However, with advances in automatic speech recognition systems being made each year, it is likely that in the future the same type of error may not occur, resulting in the inability to be able to retrieve the corresponding data file at that later date. With regard to the second problem, this is particularly significant in video data applications, since users are likely to use names and places (which may not be in the speech recognition dictionary) as input query terms. In place of these names, the automatic speech recognition system will typically replace the out of vocabulary words with a phonetically similar word or words within the vocabulary, often corrupting nearby decodings. This can also result in the failure to retrieve the required data file upon subsequent request.

In contrast, with the proposed phoneme and word lattice annotation data, a quick and efficient search using the word data in the database 29 can be carried out and, if

this fails to provide the required data file, then a further search using the more robust phoneme data can be performed. The phoneme and word lattice is an acyclic directed graph with a single entry point and a single exit point. It represents different parses of the audio stream within the data file. It is not simply a sequence of words with alternatives since each word does not have to be replaced by a single alternative, one word can be substituted for two or more words or phonemes, and the whole structure can form a substitution for one or more words or phonemes. Therefore, the density of data within the phoneme and word lattice essentially remains linear throughout the audio data, rather than growing exponentially as in the case of the N-Best technique discussed above. As those skilled in the art of speech recognition will realise, the use of phoneme data is more robust, because phonemes are dictionary independent and allow the system to cope with out of vocabulary words, such as names, places, foreign words etc. The use of phoneme data is also capable of making the system future proof, since it allows data files which are placed into the database to be retrieved even when the words were not understood by the original automatic speech recognition system.

The way in which this phoneme and word lattice annotation data can be generated for a video data file will now be described with reference to Figure 3. As shown, the video data file 31 comprises video data 31-1, which defines the sequence of images forming the video sequence and audio data 31-2, which defines the audio which is associated with the video sequence. As is well known, the audio data 31-2 is time synchronised with the video data 31-1 so that, in use, both the video and audio data are supplied to the user at the same time.



As shown in Figure 3, in this embodiment, the audio data 31-2 is input to an automatic speech recognition unit 33, which is operable to generate a phoneme lattice corresponding to the stream of audio data 31-2. Such an automatic speech recognition unit 33 is commonly available in the art and will not be described in further detail. The reader is referred to, for example, the book entitled 'Fundamentals of Speech Recognition' by Lawrence Rabiner and Bing-Hwang Juang and, in particular, to pages 42 to 50 thereof, for further information on this type of speech recognition system.

Figure 4a illustrates the form of the phoneme lattice data output by the speech recognition unit 33, for the input audio corresponding to the phrase '...tell me about Jason...'. As shown, the automatic speech recognition unit 33 identifies a number of different possible phoneme strings which correspond to this input audio utterance. For example, the speech recognition system considers that the first phoneme in the audio string is either a /t/ or a /d/. As is well known in the art of speech recognition, these different possibilities can have their own weighting which is generated by the speech recognition unit 33 and is indicative of the confidence of the speech recognition unit's output. For example, the phoneme /t/ may be given a weighting of 0.9 and the phoneme /d/ may be given a weighting of 0.1, indicating that the speech recognition system is fairly confident that the corresponding portion of audio represents the phoneme /t/, but that it still may be the phoneme /d/. In this embodiment, however, this weighting of the phonemes is not performed.

As shown in Figure 3, the phoneme lattice data 35 output by the automatic speech recognition unit 33 is input to

a word decoder 37 which is operable to identify possible words within the phoneme lattice data 35. In this embodiment, the words identified by the word decoder 37 are incorporated into the phoneme lattice data structure.

5 For example, for the phoneme lattice shown in Figure 4a, the word decoder 37 identifies the words 'tell', 'dell', 'term', 'me', 'a', 'boat', 'about', 'chase' and 'sun'. As shown in Figure 4b, these identified words are added to the phoneme lattice data structure output by the

10 speech recognition unit 33, to generate a phoneme and word lattice data structure which forms the annotation data 31-3. This annotation data 31-3 is then combined with the video data file 31 to generate an augmented video data file 31' which is then stored in the database

15 29. As those skilled in the art will appreciate, in a similar way to the way in which the audio data 31-2 is time synchronised with the video data 31-1, the annotation data 31-3 is also time synchronised and associated with the corresponding video data 31-1 and

20 audio data 31-2, so that a desired portion of the video and audio data can be retrieved by searching for and locating the corresponding portion of the annotation data 31-3.

25 In this embodiment, the annotation data 31-3 stored in the database 29 has the following general form:

#### HEADER

- time of start
- flag if word if phoneme if mixed
- 30 - time index associating the location of blocks of annotation data within memory to a given time point.
- word set used (i.e. the dictionary)
- phoneme set used
- 35 - the language to which the vocabulary

11

pertains

Block(i)  $i = 0, 1, 2, \dots$ node  $N_j$   $j = 0, 1, 2, \dots$ 

- time offset of node from start of block

- phoneme links (k)  $k = 0, 1, 2, \dots$ 

offset to node  $N_j = N_k - N_j$  ( $N_k$  is node to which link K extends) or if  $N_k$  is in block (i+1) offset to node  $N_j = N_k + N_b - N_j$  (where  $N_b$  is the number of nodes in block (i))

phoneme associated with link (k)

- word links (l)  $l = 0, 1, 2, \dots$ 

offset to node  $N_j = N_l - N_j$  ( $N_j$  is node to which link l extends) or if  $N_k$  is in block (i+1) offset to node  $N_j = N_k + N_b - N_j$  (where  $N_b$  is the number of nodes in block (i))

word associated with link (l)

The time of start data in the header can identify the time and date of transmission of the data. For example, if the video file is a news broadcast, then the time of start may include the exact time of the broadcast and the date on which it was broadcast.

The flag identifying if the annotation data is word annotation data, phoneme annotation data or if it is mixed is provided since not all the data files within the database will include the combined phoneme and word lattice annotation data discussed above, and in this case, a different search strategy would be used to search this annotation data.

In this embodiment, the annotation data is divided into blocks in order to allow the search to jump into the

middle of the annotation data for a given audio data stream. The header therefore includes a time index which associates the location of the blocks of annotation data within the memory to a given time offset between the time of start and the time corresponding to the beginning of the block.

The header also includes data defining the word set used (i.e. the dictionary), the phoneme set used and the language to which the vocabulary pertains. The header may also include details of the automatic speech recognition system used to generate the annotation data and any appropriate settings thereof which were used during the generation of the annotation data.

The blocks of annotation data then follow the header and identify, for each node in the block, the time offset of the node from the start of the block, the phoneme links which connect that node to other nodes by phonemes and word links which connect that node to other nodes by words. Each phoneme link and word link identifies the phoneme or word which is associated with the link. They also identify the offset to the current node. For example, if node  $N_{50}$  is linked to node  $N_{55}$  by a phoneme link, then the offset to node  $N_{50}$  is 5. As those skilled in the art will appreciate, using an offset indication like this allows the division of the continuous annotation data into separate blocks.

In an embodiment where an automatic speech recognition unit outputs weightings indicative of the confidence of the speech recognition units output, these weightings or confidence scores would also be included within the data structure. In particular, a confidence score would be provided for each node which is indicative of the

confidence of arriving at the node and each of the phoneme and word links would include a transition score depending upon the weighting given to the corresponding phoneme or word. These weightings would then be used to control the search and retrieval of the data files by discarding those matches which have a low confidence score.

#### DATA FILE RETRIEVAL

Figure 5 is a block diagram illustrating the form of a user terminal 59 which can be used to retrieve the annotated data files from the database 29. This user terminal 59 may be, for example, a personal computer, hand held device or the like. As shown, in this embodiment, the user terminal 59 comprises the database 29 of annotated data files, an automatic speech recognition unit 51, a search engine 53, a control unit 55 and a display 57. In operation, the automatic speech recognition unit 51 is operable to process an input voice query from the user 39 received via the microphone 7 and the input line 61 and to generate therefrom corresponding phoneme and word data. This data may also take the form of a phoneme and word lattice, but this is not essential. This phoneme and word data is then input to the control unit 55 which is operable to initiate an appropriate search of the database 29 using the search engine 53. The results of the search, generated by the search engine 53, are then transmitted back to the control unit 55 which analyses the search results and generates and displays appropriate display data to the user via the display 57.

Figures 6a and 6b are flow diagrams which illustrate the way in which the user terminal 59 operates in this embodiment. In step s1, the user terminal 59 is in an

idle state and awaits an input query from the user 39. Upon receipt of an input query, the phoneme and word data for the input query is generated in step s3 by the automatic speech recognition unit 51. The control unit 55 then instructs the search engine 53, in step s5, to perform a search in the database 29 using the word data generated for the input query. The word search employed in this embodiment is the same as is currently being used in the art for typed keyword searches, and will not be described in more detail here. If in step s7, the control unit 55 identifies from the search results, that a match for the user's input query has been found, then it outputs the search results to the user via the display 57.

In this embodiment, the user terminal 59 then allows the user to consider the search results and awaits the user's confirmation as to whether or not the results correspond to the information the user requires. If they are, then the processing proceeds from step s11 to the end of the processing and the user terminal 59 returns to its idle state and awaits the next input query. If, however, the user indicates (by, for example, inputting an appropriate voice command) that the search results do not correspond to the desired information, then the processing proceeds from step s11 to step s13, where the search engine 53 performs a phoneme search of the database 29. However, in this embodiment, the phoneme search performed in step s13 is not of the whole database 29, since this could take several hours depending on the size of the database 29.

Instead, the phoneme search performed in step s13 uses the results of the word search performed in step s5 to identify one or more portions within the database which

may correspond to the user's input query. The way in which the phoneme search performed in step s13 is performed in this embodiment, will be described in more detail later. After the phoneme search has been performed, the control unit 55 identifies, in step s15, if a match has been found. If a match has been found, then the processing proceeds to step s17 where the control unit 55 causes the search results to be displayed to the user on the display 57. Again, the system then awaits the user's confirmation as to whether or not the search results correspond to the desired information. If the results are correct, then the processing passes from step s19 to the end and the user terminal 59 returns to its idle state and awaits the next input query. If however, the user indicates that the search results do not correspond to the desired information, then the processing proceeds from step s19 to step s21, where the control unit 55 is operable to ask the user, via the display 57, whether or not a phoneme search should be performed of the whole database 29. If in response to this query, the user indicates that such a search should be performed, then the processing proceeds to step s23 where the search engine performs a phoneme search of the entire database 29.

On completion of this search, the control unit 55 identifies, in step s25, whether or not a match for the user's input query has been found. If a match is found, then the processing proceeds to step s27 where the control unit 55 causes the search results to be displayed to the user on the display 57. If the search results are correct, then the processing proceeds from step s29 to the end of the processing and the user terminal 59 returns to its idle state and awaits the next input query. If, on the other hand, the user indicates that the

search results still do not correspond to the desired information, then the processing passes to step s31 where the control unit 55 queries the user, via the display 57, whether or not the user wishes to redefine or amend the search query. If the user does wish to redefine or amend the search query, then the processing returns to step s3 where the user's subsequent input query is processed in a similar manner. If the search is not to be redefined or amended, then the search results and the user's initial input query are discarded and the user terminal 59 returns to its idle state and awaits the next input query.

#### PHONEME SEARCH

As mentioned above, in steps s13 and s23, the search engine 53 compares the phoneme data of the input query with the phoneme data in the phoneme and word lattice annotation data stored in the database 29. Various techniques can be used including standard pattern matching techniques such as dynamic programming, to carry out this comparison. In this embodiment, a technique which we refer to as M-GRAMS is used. This technique was proposed by Ng, K. and Zue, V.W. and is discussed in, for example, the paper entitled "Subword unit representations for spoken document retrieval" published in the proceedings of Eurospeech 1997.

The problem with searching for individual phonemes is that there will be many occurrences of each phoneme within the database. Therefore, an individual phoneme on its own does not provide enough discriminability to be able to match the phoneme string of the input query with the phoneme strings within the database. Syllable sized units, however, are likely to provide more discriminability, although they are not easy to identify.



The M-GRAM technique presents a suitable compromise between these two possibilities and takes overlapping fixed size fragments, or M-GRAMS, of the phoneme string to provide a set of features. This is illustrated in Figure 8, which shows part of an input phoneme string having phonemes a, b, c, d, e, and f, which are split into four M-GRAMS (a, b, c), (b, c, d), (c, d, e) and (d, e, f). In this illustration, each of the four M-GRAMS comprises a sequence of three phonemes which is unique and represents a unique feature ( $f_i$ ) which can be found within the input phoneme string.

Therefore, referring to Figure 7, the first step s51 in performing the phoneme search in step s13 shown in Figure 6, is to identify all the different M-GRAMS which are in the input phoneme data and their frequency of occurrence. Then, in step s53, the search engine 53 determines the frequency of occurrence of the identified M-GRAMS in the selected portion of the database (identified from the word search performed in step s5 in Figure 6). To illustrate this, for a given portion of the database and for the example M-GRAMS illustrated in Figure 8, this yields the following table of information:

M-GRAM (feature ( $f_i$ ))	Input phoneme string frequency of occurrence ( $q$ )	Phoneme string of selected portion of database ( $a$ )
$M_1$	1	0
$M_2$	2	2
$M_3$	3	2
$M_4$	1	1

Next, in step s55, the search engine 53 calculates a similarity score representing a similarity between the phoneme string of the input query and the phoneme string of the selected portion from the database. In this embodiment, this similarity score is determined using a cosine measure using the frequencies of occurrence of the identified M-GRAMS in the input query and in the selected portion of the database as vectors. The philosophy behind this technique is that if the input phoneme string is similar to the selected portion of the database phoneme string, then the frequency of occurrence of the M-GRAM features will be similar for the two phoneme strings. Therefore, if the frequencies of occurrence of the M-GRAMS are considered to be vectors (i.e. considering the second and third columns in the above table as vectors), then if there is a similarity between the input phoneme string and the selected portion of the database, then the angle between these vectors should be small. This is illustrated in Figure 9 for two-dimensional vectors  $\underline{a}$  and  $\underline{g}$ , with the angle between the vectors given as  $\theta$ . In the example shown in Figure 8, the vectors  $\underline{a}$  and  $\underline{g}$  will be four dimensional vectors and the similarity score can be calculated from:

$$SCORE = \cos \theta = \frac{\underline{a} \cdot \underline{g}}{|\underline{a}| |\underline{g}|} \quad (1)$$

This score is then associated with the current selected portion of the database and stored until the end of the search. In some applications, the vectors used in the calculation of the cosine measure will be the logarithm of these frequencies of occurrences, rather than the frequencies of occurrences themselves.

The processing then proceeds to step s57 where the search engine 53 identifies whether or not there are any more selected portions of phoneme strings from the database 29. If there are, then the processing returns to step s53 where a similar procedure is followed to identify the score for this portion of the database. If there are no more selected portions, then the searching ends and the processing returns to step s15 shown in Figure 6, where the control unit considers the scores generated by the search engine 53 and identifies whether or not there is a match by, for example, comparing the calculated scores with a predetermined threshold value.

As those skilled in the art will appreciate, a similar matching operation will be performed in step s23 shown in Figure 6. However, since the entire database is being searched, this search is carried out by searching each of the blocks discussed above in turn.

#### ALTERNATIVE EMBODIMENTS

As those skilled in the art will appreciate, this type of phonetic and word annotation of data files in a database provides a convenient and powerful way to allow a user to search the database by voice. In the illustrated embodiment, a single audio data stream was annotated and stored in the database for subsequent retrieval by the user. As those skilled in the art will appreciate, when the input data file corresponds to a video data file, the audio data within the data file will usually include audio data for different speakers. Instead of generating a single stream of annotation data for the audio data, separate phoneme and word lattice annotation data can be generated for the audio data of each speaker. This may be achieved by identifying, from the pitch or from another distinguishing feature of the speech signals, the

audio data which corresponds to each of the speakers and then by annotating the different speaker's audio separately. This may also be achieved if the audio data was recorded in stereo or if an array of microphones were  
5 used in generating the audio data, since it is then possible to process the audio data to extract the data for each speaker.

Figure 10 illustrates the form of the annotation data in  
10 such an embodiment, where a first speaker utters the words "... this so" and the second speaker replies "yes". As illustrated, the annotation data for the different speakers' audio data are time synchronised, relative to each other, so that the annotation data is still time  
15 synchronised to the video and audio data within the data file. In such an embodiment, the header information in the data structure should preferably include a list of the different speakers within the annotation data and, for each speaker, data defining that speaker's language,  
20 accent, dialect and phonetic set, and each block should identify those speakers that are active in the block.

In the above embodiments, a speech recognition system was used to generate the annotation data for annotating a  
25 data file in the database. As those skilled in the art will appreciate, other techniques can be used to generate this annotation data. For example, a human operator can listen to the audio data and generate a phonetic and word transcription to thereby manually generate the annotation  
30 data.

In the above embodiments, the annotation data was generated from audio stored in the data file itself. As  
35 those skilled in the art will appreciate, other techniques can be used to input the annotation data.

Figure 11 illustrates the form of a user terminal 59 which allows a user to input voice annotation data via the microphone 7 for annotating a data file 91 which is to be stored in the database 29. In this embodiment, the data file 91 comprises a two dimensional image generated by, for example, a camera. The user terminal 59 allows the user 39 to annotate the 2D image with an appropriate annotation which can be used subsequently for retrieving the 2D image from the database 29. In this embodiment, the input voice annotation signal is converted, by the automatic speech recognition unit 51, into phoneme and word lattice annotation data which is passed to the control unit 55. In response to the user's input, the control unit 55 retrieves the appropriate 2D file from the database 29 and appends the phoneme and word annotation data to the data file 91. The augmented data file is then returned to the database 29. During this annotating step, the control unit 55 is operable to display the 2D image on the display 57 so that the user can ensure that the annotation data is associated with the correct data file 91.

The automatic speech recognition unit 51 generates the phoneme and word lattice annotation data by (i) generating a phoneme lattice for the input utterance; (ii) then identifying words within the phoneme lattice; and (iii) finally by combining the two. Figure 12 illustrates the form of the phoneme and word lattice annotation data generated for the input utterance "picture of the Taj-Mahal". As shown, the automatic speech recognition unit identifies a number of different possible phoneme strings which correspond to this input utterance. As shown in Figure 12, the words which the automatic speech recognition unit 51 identifies within the phoneme lattice are incorporated into the phoneme

lattice data structure. As shown, for the example phrase, the automatic speech recognition unit 51 identifies the words "picture", "of", "off", "the", "other", "ta", "tar", "jam", "ah", "hal", "ha" and "al". The control unit 55 is then operable to add this annotation data to the 2D image data file 91 which is then stored in a database 29.

As those skilled in the art will appreciate, this embodiment can be used to annotate any kind of image such as x-rays of patients, 3D videos of, for example, NMR scans, ultrasound scans etc. It can also be used to annotate one-dimensional data, such as audio data or seismic data.

In the above embodiment, a data file was annotated from a voiced annotation. As those skilled in the art will appreciate, other techniques can be used to input the annotation. For example, Figure 13 illustrates the form of a user terminal 59 which allows a user to input typed annotation data via the keyboard 3 for annotating a data file 91 which is to be stored in a database 29. In this embodiment, the typed input is converted, by the phonetic transcription unit 75, into the phoneme and word lattice annotation data (using an internal phonetic dictionary (not shown)) which is passed to the control unit 55. In response to the user's input, the control unit 55 retrieves the appropriate 2D file from the database 29 and appends the phoneme and word annotation data to the data file 91. The augmented data file is then returned to the database 29. During this annotating step, the control unit 55 is operable to display the 2D image on the display 57 so that the user can ensure that the annotation data is associated with the correct data file 91.

Figure 14 illustrates the form of the phoneme and word lattice annotation data generated for the input utterance "picture of the Taj-Mahal". As shown in Figure 2, the phoneme and word lattice is an acyclic directed graph with a single entry point and a single exit point. It represents different parses of the user's input. As shown, the phonetic transcription unit 75 identifies a number of different possible phoneme strings which correspond to the typed input.

Figure 15 is a block diagram illustrating a document annotation system. In particular, as shown in Figure 15, a text document 101 is converted into an image data file by a document scanner 103. The image data file is then passed to an optical character recognition (OCR) unit 105 which converts the image data of the document 101 into electronic text. This electronic text is then supplied to a phonetic transcription unit 107 which is operable to generate phoneme and word annotation data 109 which is then appended to the image data output by the scanner 103 to form a data file 111. As shown, the data file 111 is then stored in the database 29 for subsequent retrieval. In this embodiment, the annotation data 109 comprises the combined phoneme and word lattice described above which allows the user to subsequently retrieve the data file 111 from the database 29 by a voice query.

Figure 16 illustrates a modification to the document annotation system shown in Figure 15. The difference between the system shown in Figure 16 and the system shown in Figure 15 is that the output of the optical character recognition unit 105 is used to generate the data file 113, rather than the image data output by the scanner 103. The rest of the system shown in Figure 16 is the same as that shown in Figure 15 and will not be

described further.

Figure 17 shows a further modification to the document annotation system shown in Figure 15. In the embodiment shown in Figure 17, the input document is received by a facsimile unit 115 rather than a scanner 103. The image data output by the facsimile unit is then processed in the same manner as the image data output by the scanner 103 shown in Figure 15, and will not be described again.

In the above embodiment, a phonetic transcription unit 107 was used for generating the annotation data for annotating the image or text data. As those skilled in the art will appreciate, other techniques can be used. For example, a human operator can manually generate this annotation data from the image of the document itself.

In the above embodiment, the database 29 and the automatic speech recognition unit were both located within the user terminal 59. As those skilled in the art will appreciate, this is not essential. Figure 18 illustrates an embodiment in which the database 29 and the search engine 53 are located in a remote server 60 and in which the user terminal 59 accesses and controls data files in the database 29 via the network interface units 67 and 69 and a data network 68 (such as the Internet). In operation, the user inputs a voice query via the microphone 7 which is converted into phoneme and word data by the automatic speech recognition unit 51. This data is then passed to the control unit which controls the transmission of this phoneme and word data over the data network 68 to the search engine 53 located within the remote server 60. The search engine 53 then carries out the search in accordance with the received phoneme and word data or controls the manipulation of the



data files (for example to control the playing, forwarding or rewinding of a video file) in accordance with the received phoneme and word data. The data retrieved from the database 29 or other data relating to the search is then transmitted back, via the data network 68, to the control unit 55 which controls the display of appropriate data on the display 57 for viewing by the user 39. In this way it is possible to retrieve and control data files in the remote server 60 without using significant computer resources in the server (since it is the user terminal 59 which converts the input speech into the phoneme and word data).

In addition to locating the database 29 and the search engine 53 in the remote server 60, it is also possible to locate the automatic speech recognition unit 51 in the remote server 60. Such an embodiment is shown in Figure 19. As shown in this embodiment, the input voice query from the user is passed via input line 61 to a speech encoding unit 73 which is operable to encode the speech for efficient transfer through the data network 68. The encoded data is then passed to the control unit 55 which transmits the data over the network 68 to the remote server 60, where it is processed by the automatic speech recognition unit 51. The phoneme and word data generated by the speech recognition unit 51 for the input query is then passed to the search engine 53 for use in searching and controlling data files in the database 29. Appropriate data retrieved by the search engine 53 is then passed, via the network interface 69 and the network 68, back to the user terminal 59. This data received back from the remote server is passed via the network interface unit 67 to the control unit 55 which generates and displays appropriate data on the display 57 for viewing by the user.

In the above embodiments, the user inputs his query by voice. Figure 20 shows an alternative embodiment in which the user inputs the query via the keyboard 3. As shown, the text input via the keyboard 3 is passed to  
5 phonetic transcription unit 75 which is operable to generate a corresponding phoneme string from the input text. This phoneme string together with the words input via the keyboard 3 are then passed to the control unit 55 which initiates a search of database using the search  
10 engine 53. The way in which this search is carried out is the same as in the first embodiment and will not, therefore, be described again. As with the other embodiments discussed above, the phonetic transcription unit 75, search engine 53 and/or the database 29 may all  
15 be located in a remote server.

In the first embodiment, the audio data from the data file 31 was passed through an automatic speech recognition unit in order to generate the phoneme  
20 annotation data. In some situations, a transcript of the audio data will be present in the data file. Such an embodiment is illustrated in Figure 21. In this embodiment, the data file 81 represents a digital video file having video data 81-1, audio data 81-2 and script  
25 data 81-3 which defines the lines for the various actors in the video film. As shown, the script data 81-3 is passed through a text to phoneme converter 83, which generates phoneme lattice data 85 using a stored dictionary which translates words into possible sequences  
30 of phonemes. This phoneme lattice data 85 is then combined with the script data 81-3 to generate the above described phoneme and word lattice annotation data 81-4. This annotation data is then added to the data file 81 to generate an augmented data file 81' which is then added  
35 to the database 29. As those skilled in the art will

appreciate, this embodiment facilitates the generation of separate phoneme and word lattice annotation data for the different speakers within the video data file, since the script data usually contains indications of who is talking. The synchronisation of the phoneme and word lattice annotation data with the video and audio data can then be achieved by performing a forced time alignment of the script data with the audio data using an automatic speech recognition system (not shown).

In the above embodiments, a phoneme (or phoneme-like) and word lattice was used to annotate a data file. As those skilled in the art of speech recognition and speech processing will realise, the word "phoneme" in the description and claims is not limited to its linguistic meaning but includes the various sub-word units that are identified and used in standard speech recognition systems.